

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
15 May 2003 (15.05.2003)

PCT

(10) International Publication Number
WO 03/040964 A2

(51) International Patent Classification⁷: **G06F 17/30**

(21) International Application Number: PCT/US02/35306

(22) International Filing Date:
4 November 2002 (04.11.2002)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
60/335,542 2 November 2001 (02.11.2001) US

(71) Applicant: **SIEMENS MEDICAL SOLUTIONS USA, INC.** [US/US]; 51 Valley Stream Parkway, Malvern, PA 19355 (US).

NICULESCU, Radu, Stefan; 2041 Wightman Street, Apt. B11, Pittsburgh, PA 15217 (US). **GOEL, Arun, Kumar**; 781 S. Middlesex Avenue, Colonia, NJ 07067 (US).

(74) Agents: **PASCHBURG, Donald, B.** et al.; Siemens Corporation, Intellectual Property Dept., 186 Wood Ave. South, Iselin, NJ 08830 (US).

(81) Designated States (*national*): CA, CN, JP.

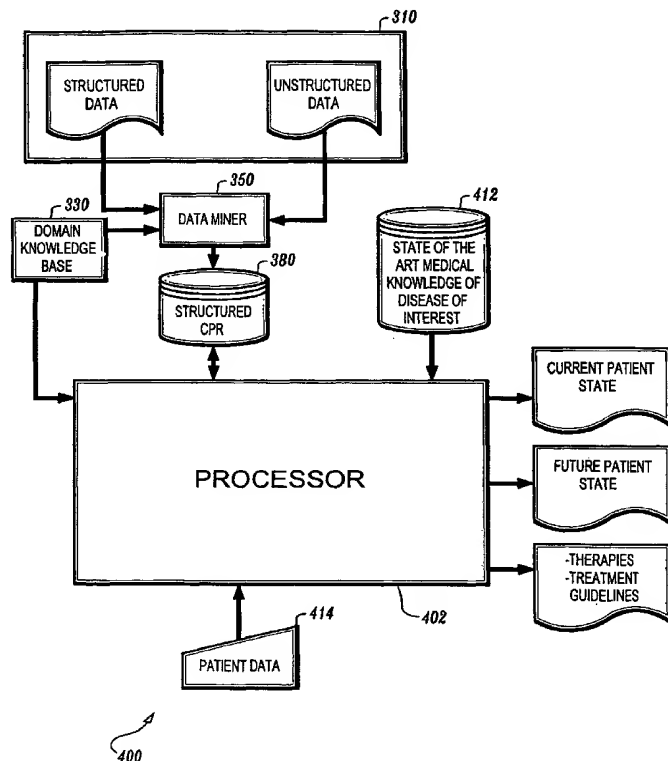
(84) Designated States (*regional*): European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR).

Published:
— *without international search report and to be republished upon receipt of that report*

(72) Inventors: **RAO, Bharat, R.**; 2060 St. Andrews Drive, Berwyn, PA 19312 (US). **SANDILYA, Sathyakama**; 28-12 Pheasant Hollow Drive, Plainsboro, NJ 08536 (US).

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: PATIENT DATA MINING FOR DIAGNOSIS AND PROJECTIONS OF PATIENT STATES



(57) Abstract: A method and system for determining patient states is provided. The method includes the steps of data mining a patient record using a domain knowledge base relating to a disease of interest (502); inputting the mined data into a model of the disease of interest (512); and determining a state of the patient based on the model (514). The system includes a data miner (350) for mining information from a patient record (310) using a domain knowledge base (330) relating to a disease of interest; and a processor (402) for creating a patient model of the disease of interest, processing the mined data in the model to determine a current state of the patient and future states for different courses of treatment, and recommending a therapy based on the determined future state.

PATIENT DATA MINING FOR DIAGNOSIS AND PROJECTIONS OF PATIENT
STATES

Cross Reference to Related Applications

This application claims the benefit of U.S. Provisional Application Serial No. 60/335,542, filed on November 2, 2001, which is incorporated by reference herein in its entirety.

Field of the Invention

The present invention relates to medical information processing systems, and, more particularly to a computerized system and method for diagnosing a current state, or condition, of a patient, projecting a future state of the patient based on various available treatment options, and recommending a course of therapy.

Background of the Invention

The major challenge facing health care providers in the present climate is to achieve a balance between a desire to reduce costs and the overriding need to maintain quality in patient care. The attempt to reduce costs without compromising quality centers around a two-fold effort to eliminate wasteful practices, and to concentrate resources on identifying those patients with the greatest likelihood of poor outcomes. By their very nature, both efforts require the use of accurate and comprehensive databases that can be extracted and analyzed

to provide a basis for intervention. Two such areas with potential for intervention are the identification of high-risk patients that would benefit from proactive approaches, e.g., by determining their future states, and the elimination wasteful practices that increase cost without a commensurate improvement in quality, or prolong length of stay, e.g., by accurately diagnosing their current state.

The problem that confronts any such effort, however, is the lack of high-quality data that can be extracted and analyzed in any meaningful or reliable way, since most hospital databases are created in text-based or other non-structured formats. Most hospitals either resort to the use of random sampling to manually review a small proportion of patient charts, or focus on relatively easily available structured information (based, for example, on DRG or ICD-9 codes) to guide their decision-making. Any truly comprehensive changes are thus left to an imperfect process, or must await a prospective data-entry system that has the capability of acting as an adequate repository of all the differing formats in which patient data are stored. At the present time managing all these different formats presents a formidable challenge in even one hospital database, let alone in different systems.

In view of the above, there exists a need for techniques to collect patient information from a variety of sources to quickly and efficiently diagnose a current state or condition of a patient and to project the future state of the patient to

help quickly identify high-risk patients, and to determine cost-effective treatments and/or therapies.

Summary of the Invention

A system and method for determining states or conditions of a patient is provided.

According to one aspect of the present invention, a method for determining patient states is provided including the steps of data mining a patient record using a domain knowledge base relating to a disease of interest; inputting the mined data into a model of the disease of interest; and determining a state of the patient based on the model.

According to another aspect of the present invention, a system for determining patient states includes a data miner for mining information from a patient record using a domain knowledge base relating to a disease of interest; and a processor for creating a patient model of the disease of interest, processing the mined data in the model to determine a current state of the patient and future states for different courses of treatment and recommending a therapy based on the estimated future disease states.

Brief Description of the Drawings

The above and other aspects, features and advantages of the present invention will become more apparent from the

following detailed description when taken in conjunction with the accompanying drawings in which:

FIG. 1 is a block diagram of a computer processing system to which the present invention may be applied according to an embodiment of the present invention;

FIG. 2 illustrates an exemplary computerized patient record (CPR); and

FIG. 3 illustrates an exemplary data mining framework for mining high-quality structured medical information;

FIG. 4 illustrates a block diagram of an exemplary diagnosis and projection system according to an embodiment of the present invention;

FIG. 5 illustrates a flow diagram for diagnosing and projecting patient states according to an embodiment of the present invention; and

FIG. 6 is a work flow diagram for diagnosing a current patient state, projecting a future patient state and suggesting therapies and treatment based on the patient states.

Description of Preferred Embodiments

To facilitate a clear understanding of the present invention, illustrative examples are provided herein which describe certain aspects of the invention. However, it is to be appreciated that these illustrations are not meant to limit

the scope of the invention, and are provided herein to illustrate certain concepts associated with the invention.

A system and method for determining states or conditions of a patient is provided. By data mining information from various sources, e.g., structured and unstructured, the present invention can gather all the information available in a patient record and use this gathered information to make a probabilistic assertions concerning prior states and a current state of a particular patient. The prior and current states of the patient can then be used in a patient model to determine future states of the patient.

It is to be understood that the present invention may be implemented in various forms of hardware, software, firmware, special purpose processors, or a combination thereof. Preferably, the present invention is implemented in software as a program tangibly embodied on a program storage device. The program may be uploaded to, and executed by, a machine comprising any suitable architecture. Preferably, the machine is implemented on a computer platform having hardware such as one or more central processing units (CPU), a random access memory (RAM), and input/output (I/O) interface(s). The computer platform also includes an operating system and microinstruction code. The various processes and functions described herein may either be part of the microinstruction code or part of the program (or combination thereof) which is executed via the operating system. In addition, various other

peripheral devices may be connected to the computer platform such as an additional data storage device and a printing device.

It is to be understood that, because some of the constituent system components and method steps depicted in the accompanying figures are preferably implemented in software, the actual connections between the system components (or the process steps) may differ depending upon the manner in which the present invention is programmed.

FIG. 1 is a block diagram of a computer processing system 100 to which the present invention may be applied according to an embodiment of the present invention. The system 100 includes at least one processor (hereinafter processor) 102 operatively coupled to other components via a system bus 104. A read-only memory (ROM) 106, a random access memory (RAM) 108, an I/O interface 110, a network interface 112, and external storage 114 are operatively coupled to the system bus 104. Various peripheral devices such as, for example, a display device, a disk storage device (e.g., a magnetic or optical disk storage device), a keyboard, and a mouse, may be operatively coupled to the system bus 104 by the I/O interface 110 or the network interface 112.

The computer system 100 may be a standalone system or be linked to a network via the network interface 112. The network interface 112 may be a hard-wired interface. However, in various exemplary embodiments, the network interface 112

can include any device suitable to transmit information to and from another device, such as a universal asynchronous receiver/transmitter (UART), a parallel digital interface, a software interface or any combination of known or later developed software and hardware. The network interface may be linked to various types of networks, including a local area network (LAN), a wide area network (WAN), an intranet, a virtual private network (VPN), and the Internet.

The external storage 114 may be implemented using a database management system (DBMS) managed by the processor 102 and residing on a memory such as a hard disk. However, it should be appreciated that the external storage 114 may be implemented on one or more additional computer systems. For example, the external storage 114 may include a data warehouse system residing on a separate computer system.

Those skilled in the art will appreciate that other alternative computing environments may be used without departing from the spirit and scope of the present invention.

Increasingly, health care providers are employing automated techniques for information storage and retrieval. The use of a computerized patient record (CPR) to maintain patient information is one such example. As shown in Fig. 2, an exemplary CPR (200) includes information that is collected over the course of a patient's treatment. This information may include, for example, computed tomography (CT) images, X-ray images, laboratory test results, doctor progress notes,

details about medical procedures, prescription drug information, radiological reports, other specialist reports, demographic information, and billing (financial) information.

A CPR typically draws from a plurality of data sources, each of which typically reflects a different aspect of a patient's care. Structured data sources, such as financial, laboratory, and pharmacy databases, generally maintain patient information in database tables. Information may also be stored in unstructured data sources, such as, for example, free text, images, and waveforms. Often, key clinical findings are only stored within physician reports, e.g., dictations.

Fig. 3 illustrates an exemplary data mining system for mining high-quality structured clinical information using data mining techniques described in "Patient Data Mining," by Rao et al., copending U.S. Patent Application Serial No. 10/____,____, (Attorney Docket No. 8706-600) filed herewith, which is incorporated by reference in its entirety. The data mining system includes a data miner (350) that mines information from a CPR (310) using domain-specific knowledge contained in a knowledge base (330). The data miner (350) includes components for extracting information from the CPR (352), combining all available evidence in a principled fashion over time (354), and drawing inferences from this combination process (356). The mined information may be stored in a structured CPR database (380). In this manner, all

information contained in a CPR, whether from a structured or unstructured source, will be stored in a structured fashion.

The extraction component (352) deals with gleaning small pieces of information from each data source regarding a patient, which are represented as probabilistic assertions about the patient at a particular time. These probabilistic assertions are called *elements*. The combination component (354) combines all the elements that refer to the same variable at the same time period to form one unified probabilistic assertion regarding that variable. These unified probabilistic assertions are called *factoids*. The inference component (356) deals with the combination of these factoids, at the same point in time and/or at different points in time, to produce a coherent and concise picture of the progression of the patient's state over time. This progression of the patient's state is called a *state sequence*.

The present invention can build an individual model of the state of a patient. The patient state is simply a collection of variables that one may care about relating to the patient. The information of interest may include a state sequence, i.e., the value of the patient state at different points in time during the patient's treatment.

Each of the above components uses detailed knowledge regarding the domain of interest, such as, for example, a disease of interest. This domain knowledge base (330) can come in two forms. It can be encoded as an input to the

system, or as programs that produce information that can be understood by the system. The part of the domain knowledge base (330) that is input to the present form of the system may also be learned from data.

As mentioned, the extraction component (352) takes information from the CPR (310) to produce probabilistic assertions (elements) about the patient that are relevant to an instant in time or time period. This process is carried out with the guidance of the domain knowledge that is contained in the domain knowledge base (330). The domain knowledge required for extraction is generally specific to each source.

Referring to FIG. 4, an exemplary diagnosis and projection system 400 according to an embodiment of the present invention is illustrated. The system 400 includes a processor 402 for extracting information from the structured CPR database 380, for creating models of diseases of interest and for processing the extracted information in a model to project a future state of a patient.

The processor is further coupled to a second database 412 including "state of the art" information relating to a disease of interest. This information may include standard procedures, established guidelines for treatments, standardized tests for assessment, etc. Additionally, the processor 402 is adapted to receive manually inputted patient data 414 which it may process and store in the structured database 380.

Each task performed by the system 400 is performed by an executable module residing either in the processor of the system 402 and/or in a memory device (e.g., RAM, ROM, external storage, etc.) of the system.

Referring to FIGS. 4 and 5, the diagnosis and projection system will be further described along with methods for diagnosing a patient's current state, for creating patient models based on a disease of interest, and for projecting a future state of a patient based on the patient's current state and the model.

First, a patient record 310 is assembled during the course of treatment of a patient over time. Additionally, a plurality of patient records for different patients (i.e., population-based data) may be assembled for a particular hospital and stored in common data storage area as the individual patient record 310. This historical data is mined using a domain knowledge base relating to a disease of interest and compiled in a structured CPR database 380 (step 502). The patient's current data is inputted into the system either manually 414 or by mining data from current tests (step 504).

A model is created to simulate a patient with similar characteristics of the patient being diagnosed. The processor 402 generates data for the model by mining data of similar patients from population-based data sources via data miner 350 using a domain knowledge base 330 of the disease of interest

(step 506). The processor 402 will then create the model of the disease of interest based on the mined data (step 510). Additionally, the processor may compile knowledge on the disease of interest from the second medical knowledge database 412 (step 508) and refine the model with this knowledge.

Once the patient model is created, all available patient data, i.e., data mined from structured and unstructured sources and/or manually input, will be entered into the model and various simulations will be run. The processor will determine a state sequence over time for the patient based on the model (step 512). The processor can further determine a future state at a particular time t , from the state sequence, to determine a preferred treatment guideline for the patient (step 514).

The development of the method according to a preferred embodiment of the present invention will now be described below in detail.

Let S be a continuous time random process taking values in Σ that represents the state of the system. Let $T = \{t_1, t_2, \dots, t_n\}$, where $t_i < t_{i+1}$, be the n "times of interest" when S has to be inferred. Let S_i refer to the sample of S at time $t_i \in T$. Let V be the set of variables that depend upon S . Let O be set of all (probabilistic) observations for all variables, $v \in V$. Let O_i be the set of all observations "assigned" to $t_i \in T$; i.e., all observations about variables, $v \in V$, that are relevant

for this time-step t_i . Similarly, let $O_i^j(v)$ be the j -th observation for variable v assigned to t_i . Let $seq = \langle S_1, S_2, \dots, S_n \rangle$ be a random variable in Σ^n ; i.e., each realization of seq is a state sequence across T . GOAL: Estimate the most likely state sequence, seq_{MAP} , (the maximum a posteriori (MAP) estimate of seq) given O :

$$seq_{MAP} = \arg \max_{seq} P[seq | O]$$

The primary focus of our interest is estimating what happened to the patient across T , the duration of interest. The estimation of the MAP state sequence can be done in two steps, the first of which is combination of observations at a fixed point in time and the second is the propagation of these inferences across time.

Each (smoothed) o_i is in the form of an a posteriori probability of a variable given the small context that it is extracted from. All observations, $O_i^j(v)$, about a variable for a single time t_i are combined into one assertion in a straightforward manner by using Bayes' theorem:

$$P[v_i | O_i^1(v_i), \dots, O_i^k(v_i)] \propto P[v_i] \cdot \prod_{j=1}^k P[O_i^j(v_i) | v_i] \propto \frac{\prod_{j=1}^k P[v_i | O_i^j(v_i)]}{P[v_i]^{k-1}}$$

At every $t_i \in T$, the relationships among S_i and V are modeled using a Bayesian Network. Because the state process is modeled as being Markov and the state as being causative (directly or

indirectly) of all the variables that we observe, we have the following equation:

$$P[seq | O] \propto P[S_0] \cdot \prod_{i=2}^n P[S_i | S_{i-1}] \cdot \prod_{i=1}^n P[O_i | S_i]$$

$$\propto \prod_{i=2}^n \frac{P[S_i | S_{i-1}]}{P[S_i]} \cdot \prod_{i=1}^n P[S_i | O_i]$$

This equation connects the *a posteriori* probability of *seq* (any sequence of samples of the state process across time) given all observations, to $P(S_i | O_i)$, the temporally local *a posteriori* probability of the state given the observations for each time instant. Essentially, we string together the temporally local Bayesian Networks by modeling each state sample, S_i , as the cause of the next sample, S_{i+1} .

The diagnosis problem is that of estimating the patient's disease state at time t_n as follows:

$$P[S_n | O] = \sum P[seq | O]$$

where the summation runs over those sequences *seq* where the final state is equal to S_n .

Further, the method will estimate (prognosing) the patient state (or any other patient variable) at a future time t_f . The following expressions are derived from the above equations to perform the prognosis for the patient:

$$P[S_f | O] = \sum P[S_f | S_n] P[S_n | O]$$

where S_f is a future state of the patient, and

$$P[V_f | O] = \sum P[V_f | S_f] P[S_f | S_n] P[S_n | O]$$

Where V_f is a future variable of the patent.

Furthermore, the method can also be used to predict the outcome of various treatment options that the patient may undergo using the same model for the patient's disease state and other variables of interest (which include the relationships between the treatment options and the outcomes thereof). The method determines $P[S_f|O, T_i]$ for each therapy option T_i and then presents this information to physicians so that they may make more informed decisions regarding the future treatment of the patient.

FIG. 6 is a work flow diagram for diagnosing a current patient state, projecting a future patient state and suggesting therapies and treatment based on the patient states.

First, a retrospective analysis 602 of a plurality of CPRs 610 is conducted via the data miner 612, which is also referred to as a REMIND (Reliable Extraction and Meaningful Inference from Non-structured data) system. The data miner 612 uses an approximate knowledge base 614 to compile structured CPRs 616. The structured CPRs 616 are used to refine the approximate knowledge base 614 to compile a refined knowledge base 618 to be used in a diagnosis phase 604.

In the diagnosis phase 604, the data miner 612 (i.e., the REMIND system) uses the refined knowledge base 618 to interact with a specific individual's CPR 620 to determine the individual's current state 622 as described above. In

addition, the system may be configured to determine based on the patient's symptoms, a disease that the patient is at risk for, and present to the physician all the information in the patient record that is relevant to the above disease. For example, if a patient comes in to the emergency room with a chest pain, the system will recognize that the patient is likely to have an acute myocardial infarction (heart attack) and present to the doctor, any information that is available regarding the patients troponin level, ECG reports etc.

Once the retrospective analysis and diagnosis phases are complete, the system and method of the present invention can recommend therapies either passively 606 or actively 608. In the passive therapy phase 606, the system will extract CPRs of similar patients 624 to compile a knowledge base 626 of patient-specific populations to determine patterns of treatments and outcomes of the similar patients. The system will assign an outcome to the future state by finding a patient similar to the patient. The system will assign probabilities to the future states by averaging outcomes of weighted outcomes of the similar patients. This knowledge base 626 will then be used to suggest treatments and therapies to the individual patient based on the most favorable outcomes.

Alternatively, therapies will be actively determined by varying potential future treatments and, in turn, projecting future patient states from the future treatments 630. The data miner 612 will use the information from the individual patient

record 620 to run various simulations with a therapy knowledge 628, which is learned from the structured database 616, therapy domain knowledge plus active feedback. Basically, the system evaluates a number of possible future treatment options (one of which is "do nothing") and projects the disease state into the future, e.g., if we put the patient on Drug 1 then what will happen".

Then, the system evaluates each of these treatments by looking at the future state of the patient. Simply, if using Drug 1 he dies with a probability of 95% but with Drug 2 he dies with a probability of 10%, the system will suggest Drug 2. The system will also consider other issues, like cost. If Drug 1 determines $P(\text{poor outcome})=84\%$ and "Do nothing" (treatment 2) is $P(\text{poor outcome}) = 85\%$, and Drug 1 costs \$5,000,000, the system might recommend against giving Drug 1. Similarly, the system will look at quality of life metrics, where if Drug 1 has severe side effects and only improves survival by 1%, it will not be recommended, or a combination of outcome, costs, quality of life, and other measures can be used to pick the best treatment.

Although illustrative embodiments of the present invention have been described herein with reference to the accompanying drawings, it is to be understood that the invention is not limited to those precise embodiments, and that various other changes and modifications may be affected

therein by one skilled in the art without departing from the scope or spirit of the invention.

What Is Claimed Is:

1. A method for determining patient states, the method comprising the steps of:

data mining a patient record using a domain knowledge base relating to a disease of interest;

inputting the mined data into a model of the disease of interest; and

determining a state of the patient based on the model.

2. The method as in claim 1, wherein the patient state is assigned a probability.

3. The method as in claim 1, further comprising the step of determining a state sequence of the patient.

4. The method as in claim 3, wherein the state sequence is assigned a probability.

5. The method as in claim 1, further comprising the step of determining a future state of the patient.

6. The method as in claim 5, wherein the future state is assigned a probability.

7. The method as in claim 1, further comprising the step of creating the model of the disease of interest by mining population-based data using the domain knowledge of the disease of interest.

8. The method as in claim 1, further comprising the step of determining a variable of the patient state.

9. The method as in claim 8, wherein the variable is assigned a probability.

10. A system for determining patient states comprising:

a data miner for mining information from a patient record using a domain knowledge base relating to a disease of interest; and

a processor for creating a patient model of the disease of interest and processing the mined data in the model to determine a state of the patient.

11. The system as in claim 10, wherein the processor assigns a probability to the patient state.

12. The system as in claim 10, wherein the processor determines a state sequence of the patient.

13. The system as in claim 12, wherein the processor assigns a probability to the state sequence.

14. The system as in claim 10, wherein the processor determines a future state of the patient.

15. The system as in claim 14, wherein the processor assigns a probability to the future state.

16. The system as in claim 10, wherein the data miner mines population-based data using the domain knowledge of the disease of interest and the processor creates the model of the disease of interest from the population-based mined data.

17. The system as in claim 14, wherein the processor assigns an outcome to the future state by finding a patient similar to the patient.

18. The system as in claim 14, wherein the processor determines a plurality of similar patients, and assigns probabilities to future states by averaging outcomes of the similar patients.

19. The system as in claim 18, wherein the processor assigns weights to the outcomes of plurality of similar patients.

20. The system as in claim 10, wherein the processor determines a variable related to the patient state.

21. The system as in claim 20, wherein the variable is assigned a probability.

22. The system as in claim 14, wherein the processor determines a plurality of outcomes by simulating a plurality of treatments based on the mined data of the patient.

23. The system as in claim 22, wherein the processor assigns probabilities to the outcomes and suggests a therapy.

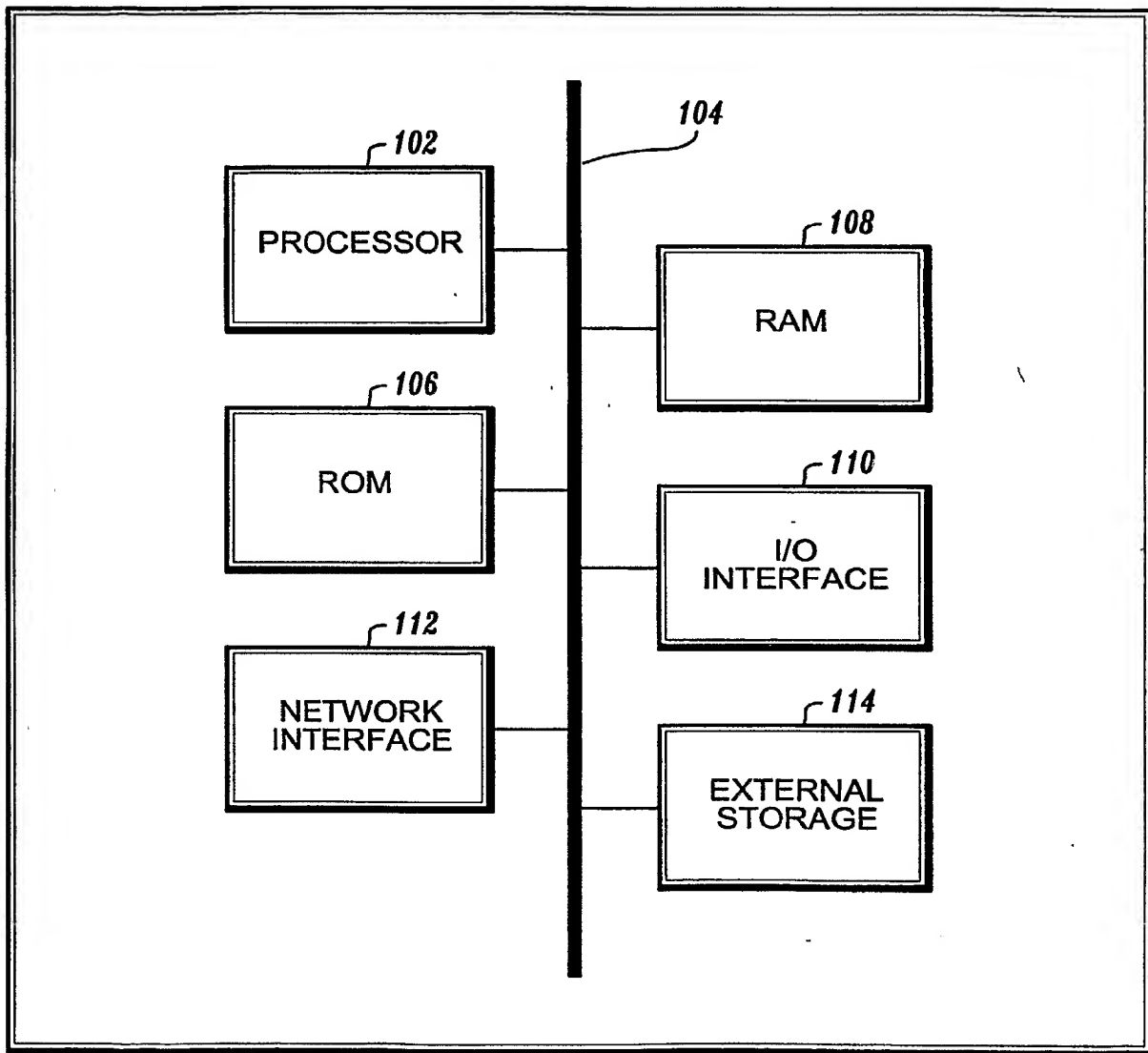
24. A program storage device readable by a machine, tangibly embodying a program of instructions executable by the machine to perform method steps for determining patient states, the method steps comprising:

data mining a patient record using a domain knowledge base relating to a disease of interest;

inputting the mined data into a model of the disease of interest; and

determining a state of the patient based on the model.

1/6

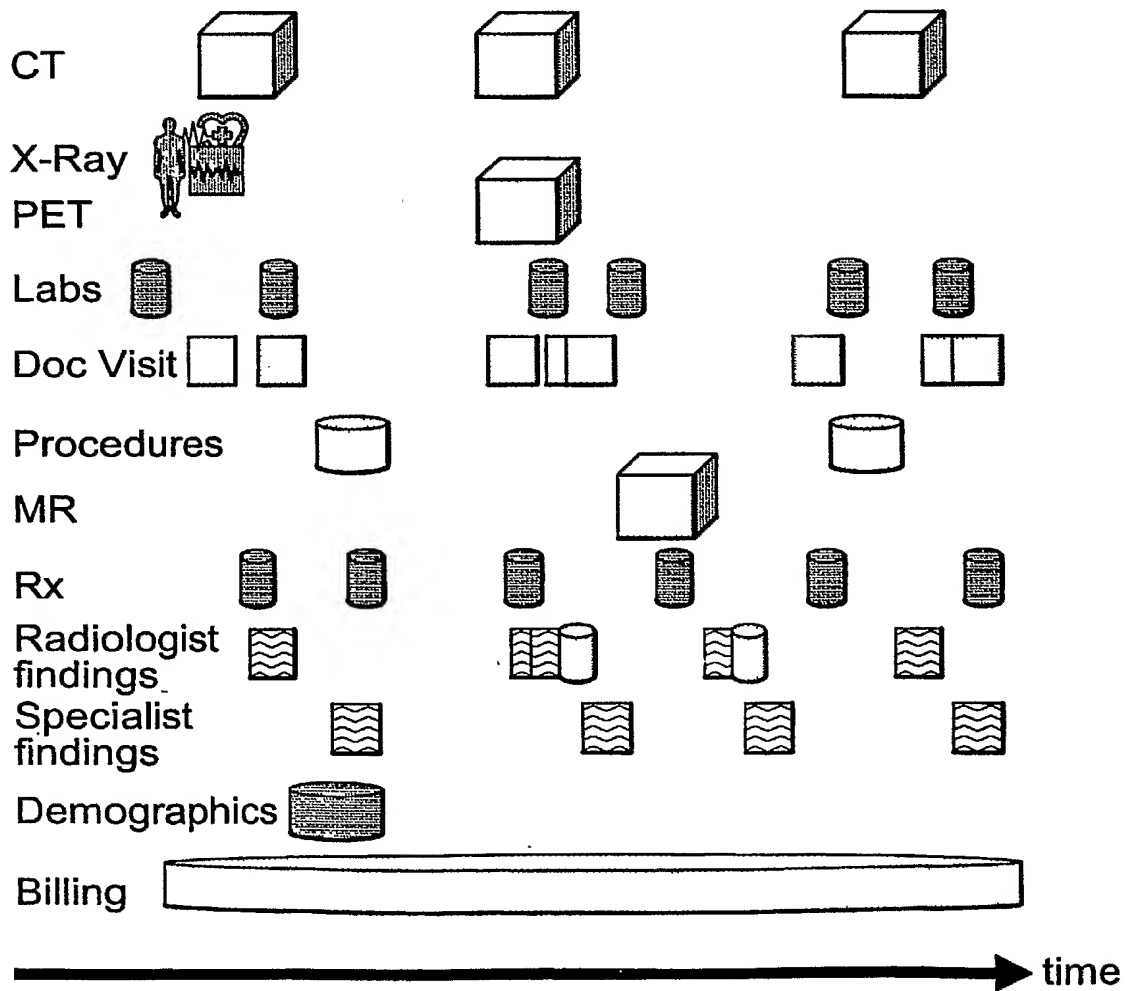


100

FIG. 1

2/6

Patient Medical Record



200

FIG. 2

3/6

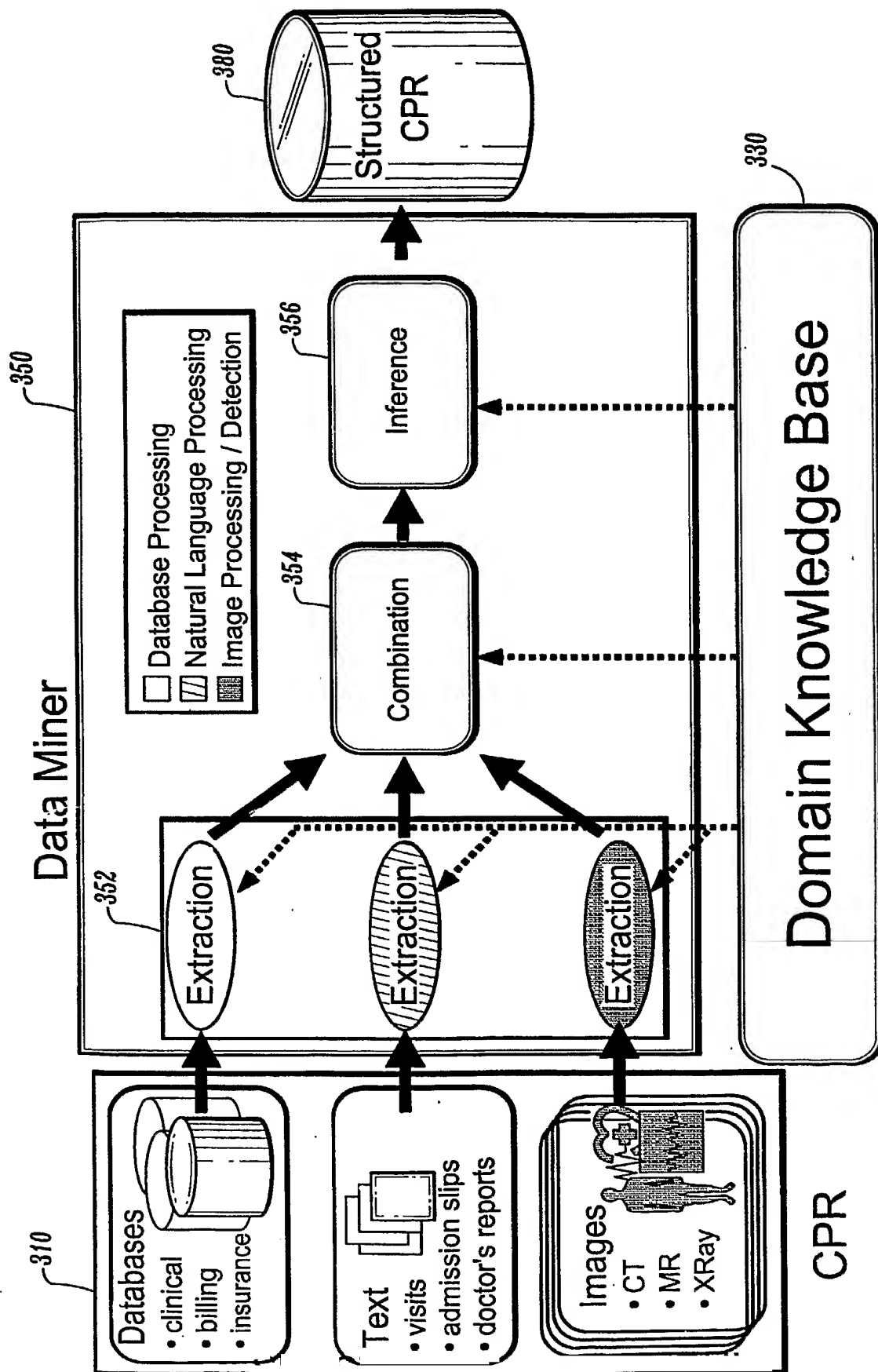
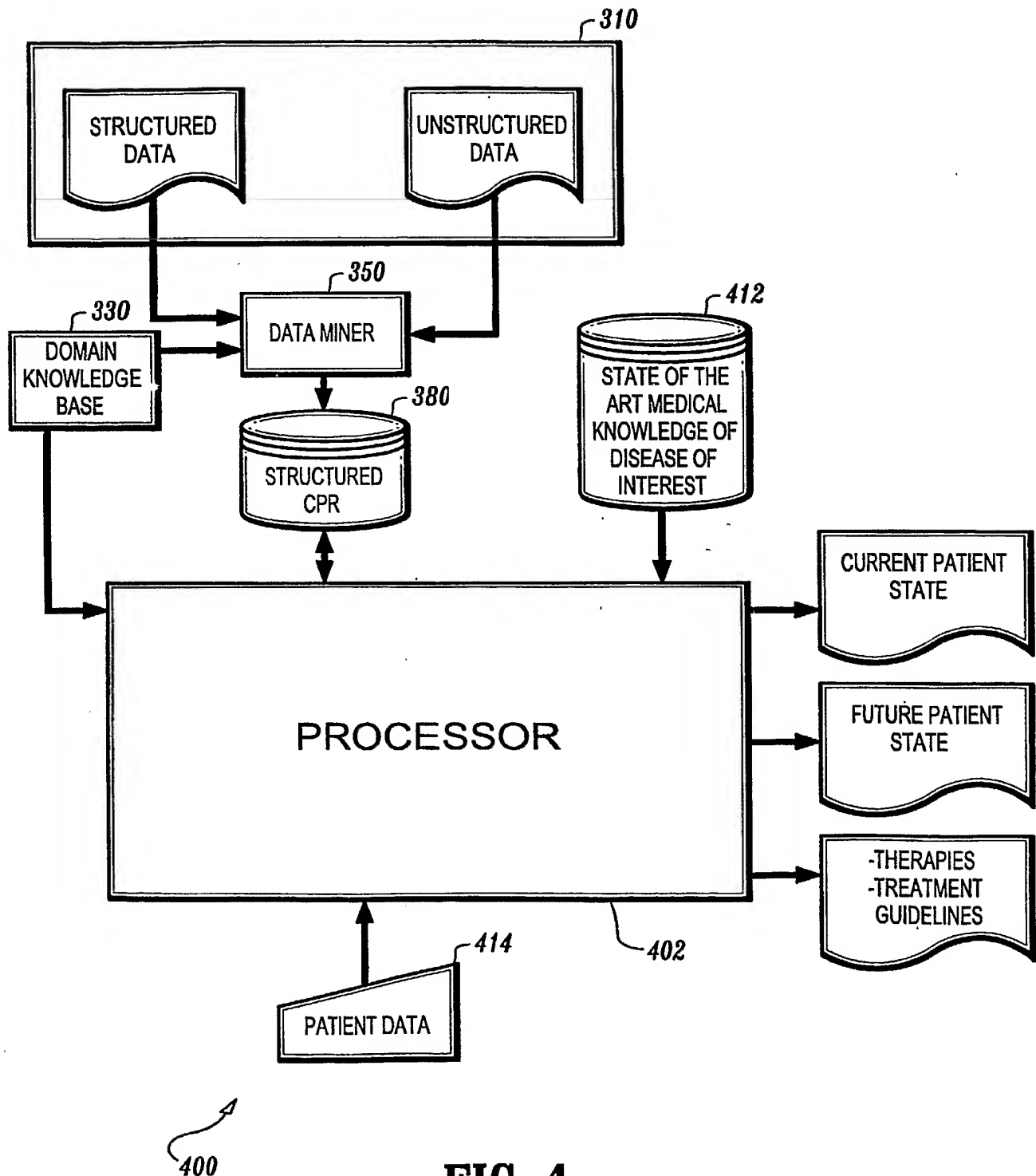
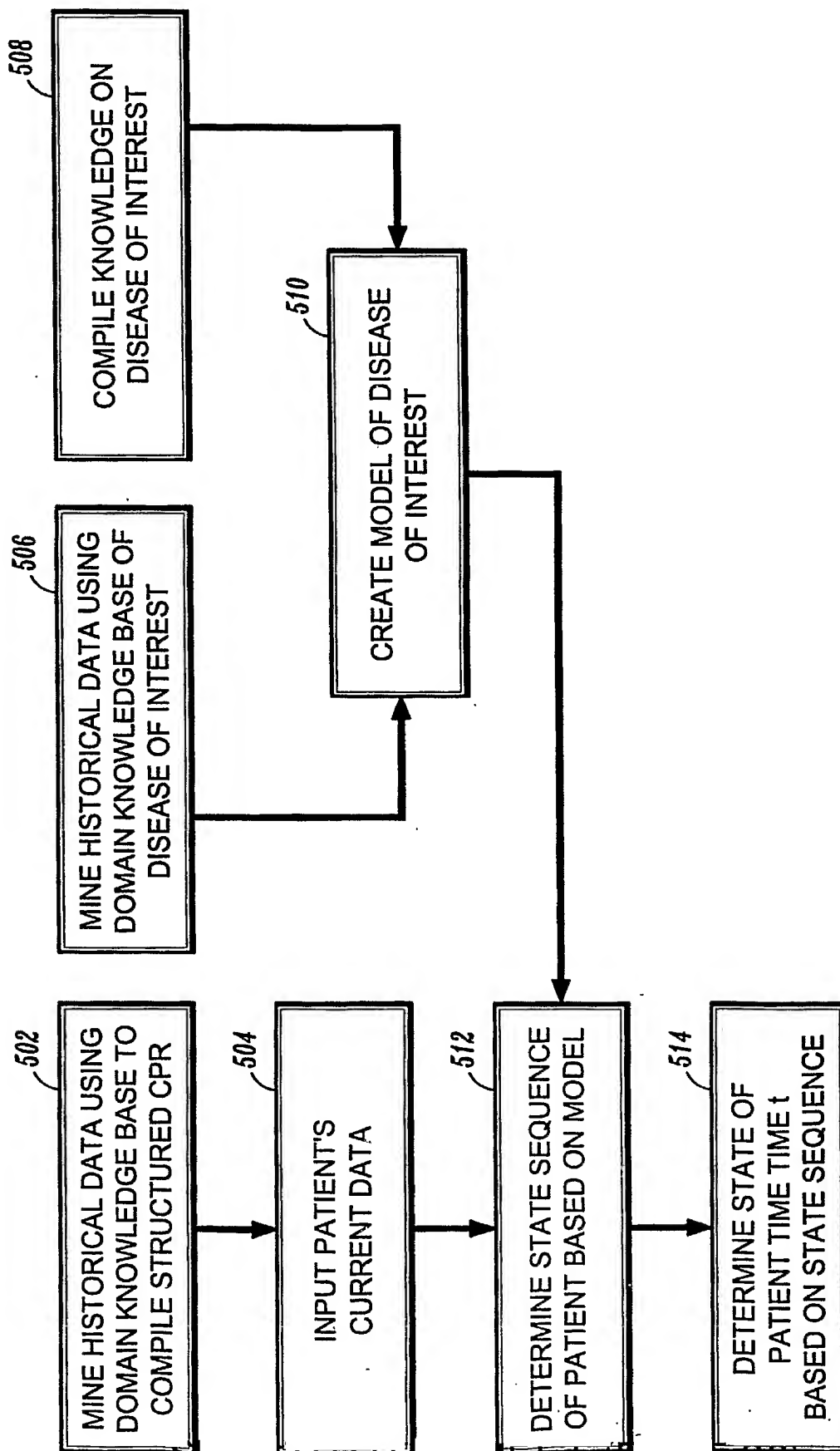


FIG. 3

4/6

**FIG. 4**

5/6

**FIG. 5**

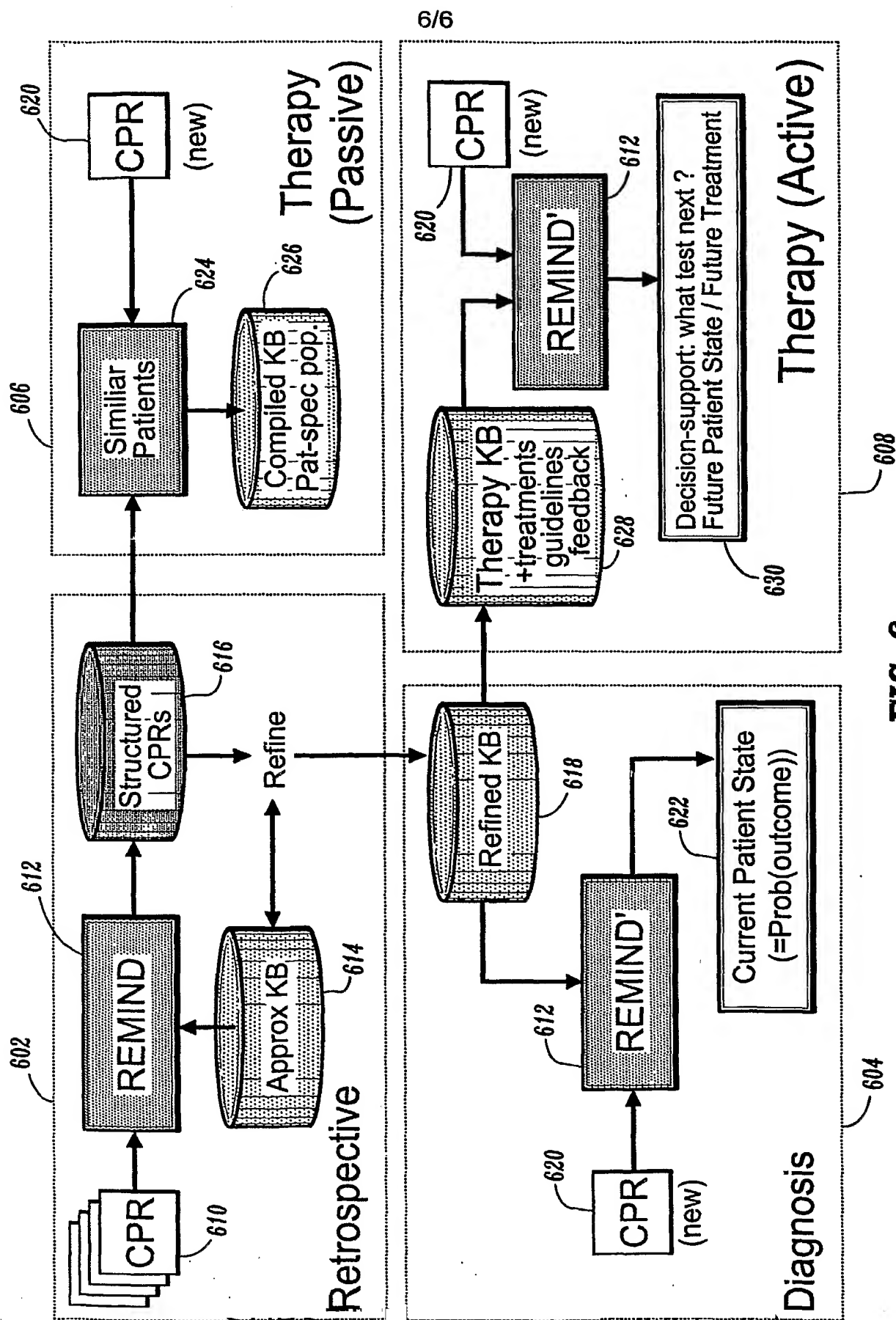


FIG. 6